

PATENT ABSTRACTS OF JAPAN

Reference 3

(11)Publication number : 11-039313

(43)Date of publication of application : 12.02.1999

(51)Int.Cl.

G06F 17/30

(21)Application number : 09-198113

(71)Applicant : NIPPON TELEGR & TELEPH CORP <NTT>

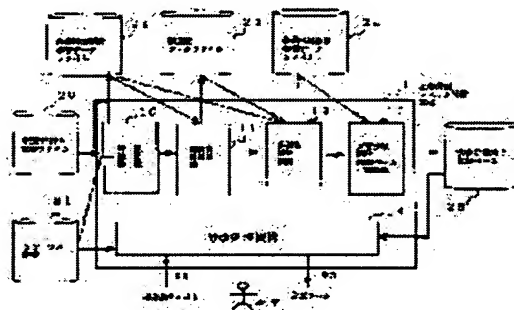
(22)Date of filing : 24.07.1997

(72)Inventor : YAMAZAKI TAKEFUMI

(54) AUTOMATIC DOCUMENT CLASSIFICATION SYSTEM, DOCUMENT CLASSIFICATION ORIENTED KNOWLEDGE BASE CREATING METHOD AND RECORD MEDIUM RECORDING ITS PROGRAM**(57)Abstract:**

PROBLEM TO BE SOLVED: To dissolve the polysemy of a semantic category and to improve precision more by selecting and utilizing only an optimum semantic category that a word has at the time of generating learning data.

SOLUTION: A characteristic extraction mechanism 10 analyzes a classification tagged document set of a file 20, examines the semantic category of a word and outputs an extracted characteristic vector and data of a classification category to which its text belongs to a learning data file 22 before polysemy dissolution. A polysemy dissolution mechanism 12 is inputted from the file 22 and an association degree data file 23 and outputs a polysemy dissolution result to a learning data file 24 after polysemy dissolution. A document classification oriented knowledge base creation mechanism 13 is inputted from the file 24 and outputs the weight between a characteristic and a classification to a document classification oriented knowledge base 25. A classification processing mechanism 14 is inputted from a thesaurus dictionary 21 and the base 25 and outputs a classification tag that corresponds to a new text inputted by a user.

**LEGAL STATUS**

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-39313

(43) 公開日 平成11年(1999) 2月12日

(51) IntCl.⁸

G 0 6 F 17/30

識別記号

F I

G 0 6 F 15/401

3 1 0 D

審査請求 未請求 請求項の数 3 O L (全 8 頁)

(21) 出願番号 特願平9-198113

(22) 出願日 平成9年(1997) 7月24日

(71) 出願人 000004226

日本電信電話株式会社

東京都新宿区西新宿三丁目19番2号

(72) 発明者 山崎 毅文

東京都新宿区西新宿三丁目19番2号 日本

電信電話株式会社内

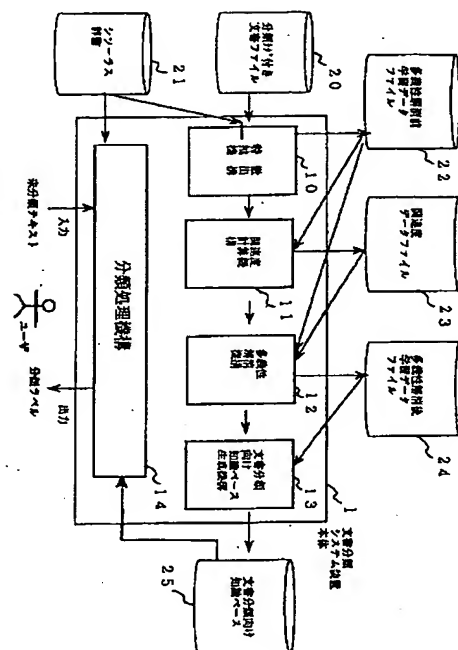
(74) 代理人 弁理士 鈴木 誠

(54) 【発明の名称】 文書自動分類システム、文書分類向け知識ベース生成方法及びそのプログラムを記録した記録媒体

(57) 【要約】

【課題】 意味カテゴリーの多義性を解消し、高精度の文書自動分類システムを提供する。

【解決手段】 分類タグ付き文書集合を入力し、シソーラス辞書を用いて、該文書の特徴付ける単語とその意味カテゴリーと該文書の分類カテゴリーからなる学習データを生成する手段、該学習データについて、分類カテゴリーと特徴間の関連度を計算する手段、該関連度に基づいて、前記学習データから不適切な特徴を除去して新たな学習データを生成する手段、該新たな学習データについて、分類カテゴリーと特徴間の重みを計算して文書分類向け知識ベースを生成する手段を設ける。



【特許請求の範囲】

【請求項 1】 シソーラス辞書を用いた文書自動分類システムにおいて、分類タグ付き文書集合とシソーラス辞書を入力とし、文書の特徴付ける単語とその意味カテゴリと当該文書の分類カテゴリから構成される学習データを生成する手段と、前記生成された学習データを入力とし、分類カテゴリと特徴間の関連度を計算する手段と、前記計算された関連度と前記学習データを入力とし、学習データから不適切な特徴を除去して新たな学習データを生成する手段と、前記生成された新たな学習データを入力とし、分類カテゴリと特徴間の重みを計算して文書分類向け知識ベースを生成する手段と、未分類の文書を入力とし、前記知識ベースを元に対応する分類カテゴリを出力する手段とを有することを特徴とする文書自動分類システム。

【請求項 2】 シソーラス辞書を用いた文書自動分類システムにおける文書分類向け知識ベースを生成する方法であって、分類タグ付き文書集合とシソーラス辞書を入力して、文書の特徴付ける単語とその意味カテゴリ及び当該文書の分類カテゴリから構成される学習データを生成し、前記学習データについて、分類カテゴリと特徴間の関連度を計算し、前記関連度に基づき、前記学習データから不適切な特徴を除去して新たな学習データを生成し、前記新たな学習データにより、分類カテゴリと特徴間の重みを計算して、文書分類向け知識ベースを生成することを特徴とする文書分類向け知識ベース生成方法。

【請求項 3】 シソーラス辞書を用いた文書自動分類システムにおける文書分類知識ベースを生成する処理のプログラムを記録したコンピュータ読み取り可能な記録媒体であって、分類タグ付き文書集合を読み取る処理と、前記分類タグ付き文書集合を解析して、当該文書の特徴付ける単語を抽出し、シソーラス辞書を参照して該単語の意味カテゴリを調べ、前記単語と意味カテゴリからなる特徴ベクトル及び当該文書の分類カテゴリから構成される学習データを生成する処理と、前記学習データについて、分類カテゴリと特徴間の関連度を計算する処理と、前記関連度に基づき、前記学習データから不適切な特徴を除去して新たな学習データを生成する処理と、前記新たな学習データにより、分類カテゴリと特徴間の重みを計算して、文書分類向け知識ベースを生成する処理と、を含むプログラムを記録したコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】**【0001】**

【発明の属する技術分野】 本発明は自然言語処理の技術分野に係わり、詳しくは、文書自動分類システム、分類

タグ付き文書集合から文書分類向け知識ベースを生成する方法及びそのプログラムを記録した記録媒体に関するものである。

【0002】

【従来の技術】 文書分類とは、テキストをその意味内容に基づいて、予め決められた分類カテゴリ（例えば、政治、スポーツ、経済等）に分類することであり、大量の電子化されたテキストが流通するようになった現在、テキストの効率的検索／利用において、重要な課題となっている。この文書分類作業の自動化を実現するのが文書自動分類システムである。このシステムは、予め、カテゴリ分類されたテキスト（ここでは、これを「分類タグ付き文書集合」と呼ぶ）を利用した適当な分類器（文書分類向け知識ベース）を作成し、その分類器に基づいて新たに入力されたテキストのカテゴリ分類を行う。

【0003】 従来の文書自動分類システムは、分類タグ付き文書集合から、分類のキーとなる単語を抜き出し、それらの単語と予め付与された分類カテゴリとの結び付きの強さである関連度を学習して、分類器を生成する。同じ分類カテゴリを持つテキストには、その分野を特徴づける単語が決まって頻出するので、その単語とその指定した分類カテゴリとの関連度は大きくなり、そうでない単語と分類カテゴリとの関連度は小さくなる。分類器は、単語、分類カテゴリ、及びこの単語と分類カテゴリとの関連度を表わすリンクから構成される。

【0004】 分類器の例を図 2 に示す。図 2 で示す通り、特徴ノードはテキストで頻出する単語、分類ノードは文書集合が持つ分類カテゴリから構成され、特徴ノードと分類ノード間のリンク上の数値が、関連度を表わす。関連度の学習方法として、カイ 2 乗計算手法等、様々な手法が考えられる。

【0005】 新たなテキストが入力されると、まず、そのテキスト上に現われる単語を抽出し、次に分類器上の単語と分類タグ間の関連度から、そのテキストの各分類カテゴリへの関連度を計算する。この値が、ある決められた閾値を超えるもの或いは最も大きい値を持つものを、そのテキストの分類カテゴリと判断する。

【0006】

【発明が解決しようとする課題】 上記従来技術によって、分類タグ付き文書集合から自動的に分類器が生成できるが、特徴量ノードを構成するものが、分類タグ付き文書集合上に出現する単語に限られるので、新たに入力される分類対象テキストが、特徴ノードにない未知単語を含む場合、そのテキストの分類精度が悪くなってしまう。この問題を解決する一つの方法として、特徴ノードの生成において、単語のみでなく、シソーラス利用によって得られる単語の意味カテゴリも合わせて、利用する手法が考えられる。意味カテゴリは、単語よりも一

一般的な概念で、通常複数の単語が共通の意味カテゴリーを持つので、この意味カテゴリーによる特徴ノードの生成で、未知単語の問題が解決できる。

【0007】しかし、一方で、シソーラス利用によって得られる各単語の意味カテゴリーの数は通常一つとは限らない。例えば、単語「路線」は、政治の分野では「方針」という意味で使われ、また、交通の分野では「道筋」という意味で使われる。このような多義語は、シソーラス上で複数の意味カテゴリーを持つ。分類器を生成する学習段階で、シソーラスを利用して特徴ノードを生成する場合に、多義語から得られる複数の意味カテゴリーを、文脈に応じて絞らずに全て利用すると、特徴ノードと分類ノードとの関連度が誤って計算されてしまい、分類精度が悪くなる可能性がある。特徴ノードにシソーラス上の意味カテゴリーを利用する場合は、意味カテゴリーを文脈に応じて正確に特定する、所謂多義解消が必要である。

【0008】本発明の目的は、意味カテゴリーの多義性を解消し、より精度の高い文書自動分類システム、その文書分類向け知識ベースの生成方法及びそのプログラムを記録した記録媒体を提供することにある。

【0009】

【課題を解決するための手段】本発明の文書自動分類システムは、分類タグ付き文書集合とシソーラス辞書を入力し、文書の特徴付ける単語とその意味カテゴリーと当該文書の分類カテゴリーから構成される学習データを生成する手段と、前記生成された学習データから分類カテゴリーと特徴間の関連度を計算する手段と、前記計算された関連度と前記学習データを入力し、学習データから不適正な意味カテゴリーを除去して新たな学習データを生成する手段と、前記生成された新たな学習データを入力とし、分類カテゴリーと特徴間の関連度を計算して文書分類向け知識ベースを生成する手段と、未分類の文書を入力とし、前記知識ベースを元に対応する分類カテゴリーを出力する手段とを有することを特徴とする。

【0010】本発明の文書分類向け知識ベース生成方法は、分類タグ付き文書集合とシソーラス辞書を入力して、文書の特徴付ける単語とその意味カテゴリー及び当該文書の分類カテゴリーから構成される学習データを生成し、前記学習データについて、分類カテゴリーと特徴間の関連度を計算し、前記関連度に基づき、前記学習データから不適切な特徴を除去して新たな学習データを生成し、前記新たな学習データにより、分類カテゴリーと特徴間の重みを計算して、文書分類向け知識ベースを生成することを特徴とする。

【0011】本発明のコンピュータ読み取り可能な記録媒体は、分類タグ付き文書集合を読み取る処理と、前記分類タグ付き文書集合を解析して、当該文書の特徴付ける単語を抽出し、シソーラス辞書を参照して該単語の意味カテゴリーを調べ、前記単語と意味カテゴリーからな

る特徴ベクトル及び当該文書の分類カテゴリーから構成される学習データを生成する処理と、前記学習データについて、分類カテゴリーと特徴間の関連度を計算する処理と、前記関連度に基づき、前記学習データから不適切な特徴を除去して新たな学習データを生成する処理と、前記新たな学習データにより、分類カテゴリーと特徴間の重みを計算して、文書分類向け知識ベースを生成する処理とを含むことを特徴とする。

【0012】意味カテゴリーの多義性解消は、「意味カテゴリーは、その対象テキストが属する分類カテゴリーに関連のある他の単語からも生成される」という性質に基づいて、テキストの分類カテゴリー情報と事前に得られる特徴と該当分類カテゴリー間の関連度の利用により行える。例えば、多義語「路線」は、あるシソーラス上で、「S交通路」「S形勢」の2つの意味カテゴリーを持つ。ここでは、意味カテゴリーを単語と区別するため、意味カテゴリーの先頭に文字Sを付与する。「S形勢」は、分類カテゴリーが「政治」で頻出する単語「動向」「非常事態」からも生成される意味カテゴリーであり、また「S交通路」は、分類カテゴリーが「交通」で頻出する単語「新幹線」「東海道」からも生成される。よって、「政治」「交通」の分類カテゴリーを持つカテゴリーを持つテキストから、「S交通路」は「交通」と強い関連があり、「S形勢」は「政治」と強い関連を持つことが解るはずである。この関連度を利用すれば、対象テキスト中に単語「路線」が現われた時、そのテキストの分類カテゴリーが「政治」であれば、意味カテゴリーとして「S形勢」を選択し、「交通」であれば、「S交通路」を選択して用いることが可能である。

【0013】このように、シソーラス利用の文書自動分類システムにおいて、最適な意味カテゴリーを選択する、多義を解消する処理プログラムを組み込むことにより、学習データ中のノイズが減少し、より精度の高い文書分類向け知識ベースを生成することが可能になる。

【0014】

【発明の実施の形態】以下、図面を用いて、本発明の一実施例を説明する。図1は、本発明の一実施例に係わる文書自動分類システムのブロック図である。文書自動分類システムは、特徴抽出機構10、関連度計算機構11、多義性解消機構12、文書分類向け知識ベース生成機構13、及び未分類の文書に分類ラベルを出力する分類処理機構14から装置本体1と、分類タグ付き文書ファイル20、シソーラス辞書21、多義性解消前の学習データファイル22、関連度データファイル23、多義性解消後の学習データファイル24、文書分類向け知識ベース25等を格納する外部記憶装置群で構成される。装置本体1は、CPU、RAM、内蔵ハードディスクなどで構成される所謂コンピュータである。

【0015】図1の構成において、特徴抽出機構10は、分類タグ付き文書ファイル20及びシソーラス辞書

21を入力として、文書を表わす特徴及びその分類カテゴリー情報を多義性解消前の学習データファイル22に出力する。関連度計算機構11は、多義性解消前の学習データファイル22を入力として、特徴と分類カテゴリー間の関連度情報を関連度データファイル23に出力する。多義性解消機構12は、多義性解消前の学習データファイル22及び関連度データファイル23を入力として、多義性解消結果を多義性解消後の学習データファイル24に出力する。文書分類向け知識ベース生成機構13は、多義性解消後の学習データファイル24を入力として、再計算後の特徴と分類カテゴリー間の関連度を文書分類向け知識ベース25に出力する。分類処理機構14は、文書分類向け知識ベース25及びユーザが入力した新たなテキストを入力として、入力されたテキストに対応する分類カテゴリーを出力する。以下に、各機構10、11、12、13、14の構成および動作を詳述する。

【0016】〈特徴抽出機構10〉特徴抽出機構10は、ファイル20の入力として与えられた分類タグ付き文書集合を解析し、品詞が名詞、固定名詞である単語を抜き出し、その後、シソーラス辞書21を参照して、それらの単語の意味カテゴリーを調べる。そして、一つのテキストから、そのテキスト中に含まれる単語と意味カテゴリーからなる特徴ベクトルを抽出し、該特徴ベクトルとそのテキストの属する分類カテゴリーからなるデータを多義性解消前の学習データファイル22に出力する。

【0017】図3に、特徴抽出機構10の一実施例の構成図を示す。本特徴抽出機構10は、形態素解析部101、単語抽出部102、意味カテゴリー付与部103からなる。

【0018】形態素解析部101は、ファイル20から入力されたテキストについて形態素解析を行い、単語抽出とその品詞付けを行う。単語抽出部102は、必要な品詞（名詞、固有名詞）の単語を選び出す。意味カテ

リー付与部103は、単語抽出部102で選ばれた単語について、利用するシソーラス辞書21を調べて、その意味カテゴリーを付与し、単語と意味カテゴリーからなる特徴ベクトルと該テキストの属する分類カテゴリーからなるデータ（多義性解消前の学習データ）をファイル22に出力する。

【0019】図4は、本特徴抽出機構10での具体的処理例を示したものである。ここで、意味カテゴリーの先頭には文字「S」を付与し、単語と区別する。例えば、単語「路線」に対し、シソーラス辞書21において「交通路」「形勢」の二つの意味カテゴリーを持つ場合、「S交通路」、「S形勢」を付加する。「S党」、「S団体」も同様である。

【0020】〈関連度計算機構11〉関連度計算機構11は、多義性解消前の学習データファイル22を入力として、特徴と分類カテゴリー間の関連度情報を関連度データファイル23に出力する。

【0021】図5に、関連度計算機構11の一実施例の構成図を示す。本関連度計算機構11は、事例数カウンターモジュール111、カイ2乗計算モジュール112、関連度格納モジュール113から構成される。

【0022】事例数カウンターモジュール111は、多義性解消前の学習データファイル22を入力として、各分類カテゴリー毎に、各特徴に対して、次の4つの値をカウントする。即ち、該当特徴が出現するテキスト集合中で、該当分類カテゴリーを持つテキスト数「 N_{r+} 」、持たないテキストの数「 N_{n+} 」、また、該当特徴が出現しないテキスト集合中で、該当分類カテゴリーを持つテキスト数「 N_{r-} 」、持たないテキストの数「 N_{n-} 」をそれぞれ求める。

【0023】カイ2乗計算モジュール112は、上記4つの値を用いて、次の式に基づいて、カイ2乗値（関連度）を計算する。但し、 N は全事例数を表す。

【0024】

【数1】

$$\chi^2 = \frac{N(N_{rt}N_{n-} - N_{rn}N_{nt})^2}{(N_{rt} + N_{rn})(N_{nt} + N_{n-})(N_{rt} + N_{nt})(N_{rn} + N_{n-})}$$

【0025】関連度格納モジュール113は、カイ2乗計算モジュール112の計算結果（各分類カテゴリー毎、各特徴毎のカイ2乗値）を関連度データファイル23に格納する。

【0026】〈多義性解消機構12〉多義性解消機構12は、多義性解消前の学習データファイル22及び関連度データファイル23を入力として、多義性解消結果を多義性解消後の学習データファイル24に出力する。

【0027】図6は、多義性解消機構12の一実施例の構成図を示す。本多義性解消機構12は、多義語選択モジュール121、関連度参照モジュール122、最適意味カテゴリー選択モジュール123から構成される。

【0028】多義語選択モジュール121は、多義性解消前の学習データファイル22を入力として、意味カテゴリーが複数付与されている多義語を探しだし、その意味カテゴリー及びそのテキストが属する分類カテゴリーを出力する。関連度参照モジュール122は、関連度データファイル23を入力として、モジュール121で得られた意味カテゴリーと分類カテゴリーから、それらに対応する関連度を出力する。最適意味カテゴリー選択モジュール123は、モジュール122で得られた関連度を元に、関連度の最も大きい値をもつものを最適意味カテゴリーとして選択し、多義性解消後の学習データファイル24に出力する。

【0029】例えば、図4の場合、多義語「路線」に対応する「S交通路」と「S形勢」の2つの意味カテゴリーのうち、「S形勢」が最適意味カテゴリーとして選択される。

【0030】〈文書分類向け知識ベース生成機構13〉文書分類向け知識ベース生成機構13は、多義性解消後の学習データファイル24を入力として、特徴と分類タグ間の重みを文書分類向け知識ベース25に出力する。実施例の一例として、ここでは、文書分類を線形分類モデルに基づいて行ない、線形モデルにおける重みの学習に誤り駆動型学習アルゴリズムを用いるとする。線形分類モデルは、特徴集合ノードと分類カテゴリーノードとから構成されており、入力された事例がその分類カテゴリーに属するか否かの判定を、入力事例のもつ特徴ノードと分類カテゴリー間の重みの合計が、決められた閾値を超えるか否かによって行う。線形分類モデルの例を図7に示す。なお、線形分類モデル、誤り駆動型学習アルゴリズムについては、例えば「N. Littlestone, "Learning quickly when irrelevant attributes abound: A new linear threshold algorithm", Machine Learning, No. 2, pp285-318, 1988.」に記述されている。

【0031】図8は、誤り駆動型学習アルゴリズムを用いた、文書分類向け知識ベース生成機構13の一実施例の構成図を示す。本文書分類向け知識ベース生成機構13は重み初期値設定モジュール131、重み合算モジュール132、正解判定モジュール133、重み更新モジュール134から構成される。

【0032】重み初期値設定モジュール131は、多義性解消後の学習データファイル24を入力として、学習データに出現する特徴と分類カテゴリー間をある値に初期化し、文書分類向け知識ベース25に出力する。重み合算モジュール132は、多義性解消後の学習データファイル24、文書分類向け知識ベース25を入力として、各入力事例毎に、事例に出現する特徴から、各分類カテゴリーに対するスコアを重みの合計として計算する。正解判定モジュール133は、前モジュールで計算されたスコアが、決められた閾値を超えるか否かを判定し、分類カテゴリーを割り当てる。重み更新モジュール134は、前モジュールが割り当てた分類カテゴリーが正解の分類カテゴリーと異なる場合のみ、特徴と分類カテゴリー間の重みを更新し、文書分類向け知識ベース25に出力する。

【0033】〈分類処理機構14〉分類処理機構14は、シソーラス辞書21、文書分類向け知識ベース25及びユーザが入力した新たなテキストを入力として、入力されたテキストに対応する分類タグを出力する。

【0034】図9は、分類処理機構14の一実施例の構成図を示す。本分類処理機構14は特徴抽出モジュール141、重み合算モジュール142、分類カテゴリー生

成モジュール143から構成される。

【0035】特徴抽出モジュール141は、ユーザが入力した新たな未分類テキストを形態素解析し、名詞、固定名詞である単語を選択し、それらの持つ意味カテゴリーを、シソーラス辞書21を参照して付与し、特徴ベクトルを生成する。本モジュール141は、前記特徴抽出機構11と基本的に同じものである。重み合算モジュール142は、モジュール141で生成された特徴ベクトルと文書分類向け知識ベース25を入力として、各分類カテゴリーに対するスコア計算を行う。分類カテゴリー生成モジュール143は、前モジュールで計算されたスコアがある決められた閾値以上である分類カテゴリーを出力する。

【0036】以上、本発明の一実施例に係わる文書自動分類システムについて説明したが、図1において、文書分類向け知識ベースの生成に関係する特徴抽出機構10、関連度計算機構11、多義性解消機構12及び文書分類向け知識ベース生成機構13の処理プログラムは一つにまとめてもよい。図10は、その処理フローを示したもので、処理1010～1013は図1の各機構10～13に対応する。処理1010～1013のプログラムは、あらかじめCD-ROM等の記録媒体に記録しておき、該プログラムをコンピュータにロードすることにより、先に説明した図1の各機構10～13と同様の処理が実現する。

【0037】

【発明の効果】以上説明したように、本発明によれば、分類タグ付き文書集合とシソーラス辞書から、分類器（文書分類向け知識ベース）の生成に利用する学習データを作成する際に、単語の持つ意味カテゴリーを最適な意味カテゴリーのみを選んで利用するので、従来手法による多義性解消をしない場合に比べて、より精度の高い分類器を生成することができる。

【図面の簡単な説明】

【図1】本発明の一実施例の文書自動分類システムの全体構成図である。

【図2】分類器の一例を示す図である。

【図3】図1の特徴抽出機構の一実施例の構成図である。

【図4】特徴抽出機構の具体的処理例を示す図である。

【図5】図1の関連度計算機構の一実施例の構成図である。

【図6】図1の多義性解消機構の一実施例の構成図である。

【図7】線形分類モデルの一例を示す図である。

【図8】図1の文書分類向け知識ベース生成機構の一実施例の構成図である。

【図9】図1の分類処理機構の一実施例の構成図である。

【図10】本発明の文書分類向け知識ベース生成方法の

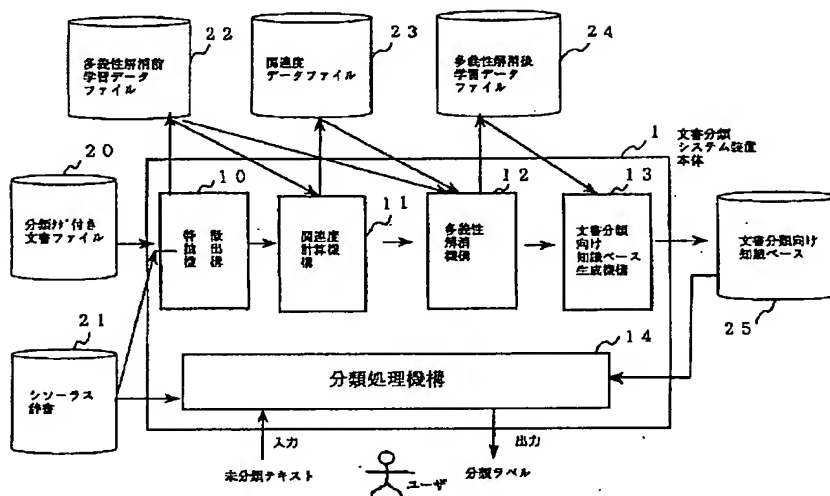
概略処理フロー図である。

【符号の説明】

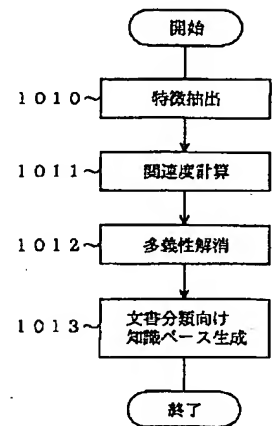
- 1 文書自動分類装置本体
- 2 外部記憶装置群
- 10 特徴抽出機構
- 11 関連度計算機構
- 12 多義性解消機構
- 13 文書分類向け知識ベース生成機構

- 14 分類処理機構
- 20 分類タグ付き文書ファイル
- 21 シソーラス辞書
- 22 多義性解消前の学習データファイル
- 23 関連度データファイル
- 24 多義性解消後の学習データファイル
- 25 文書分類向け知識ベース

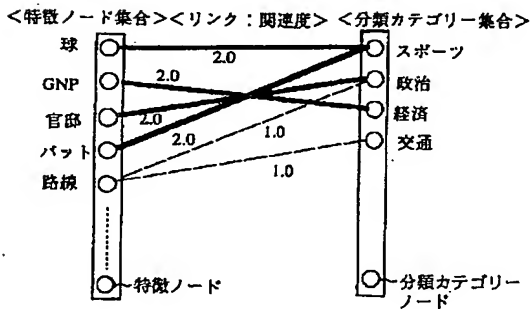
【図 1】



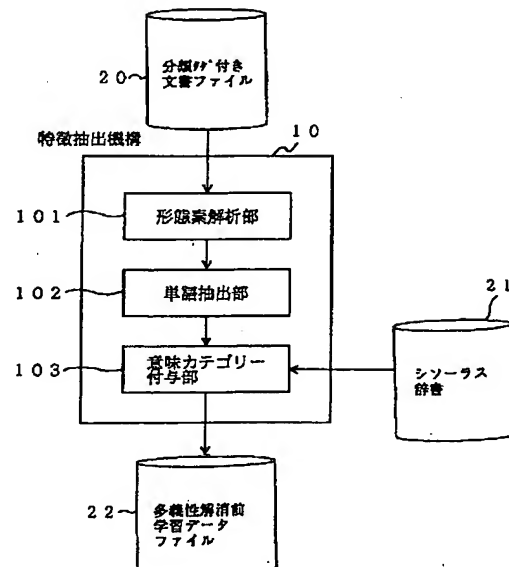
【図 10】



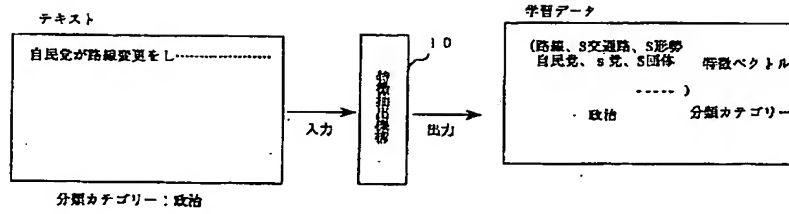
【図 2】



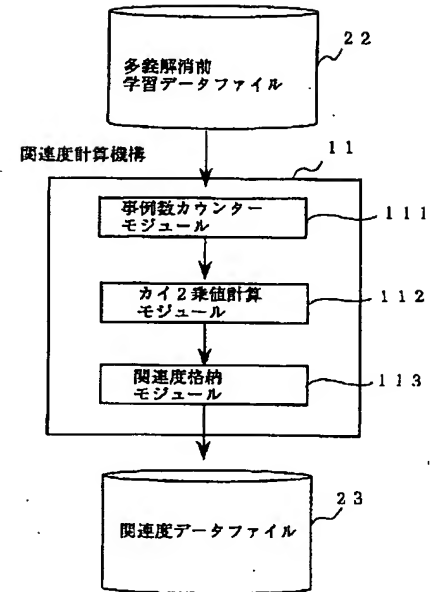
【図 3】



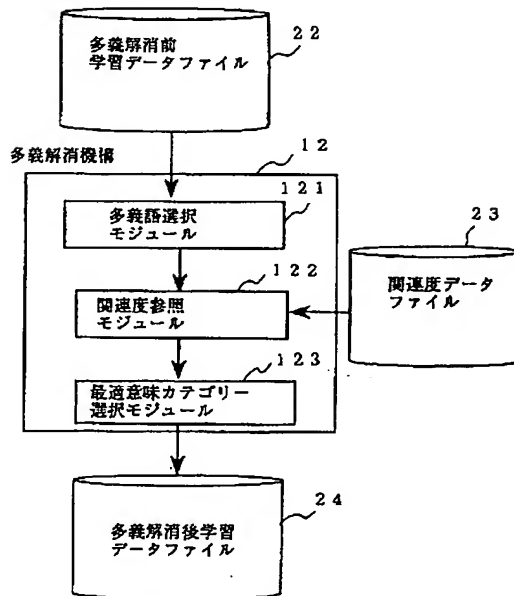
【図4】



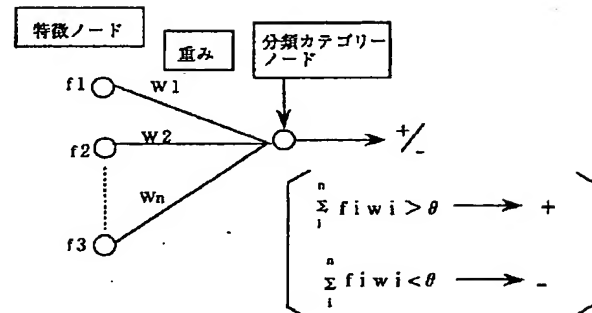
【図5】



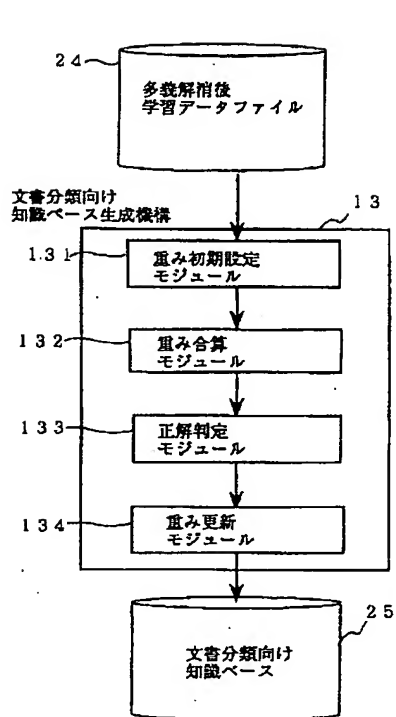
【図6】



【図7】



【図 8】



【図 9】

